

A Relation-Based Page Rank Algorithm for Semantic Web Search Engines

*Original*

A Relation-Based Page Rank Algorithm for Semantic Web Search Engines / Lamberti, Fabrizio; Sanna, Andrea; Demartini, Claudio Giovanni. - In: IEEE TRANSACTIONS ON KNOWLEDGE AND DATA ENGINEERING. - ISSN 1041-4347. - STAMPA. - 21:1(2009), pp. 123-136. [10.1109/TKDE.2008.113]

*Availability:*

This version is available at: 11583/1801015 since:

*Publisher:*

*Published*

DOI:10.1109/TKDE.2008.113

*Terms of use:*

openAccess

This article is made available under terms and conditions as specified in the corresponding bibliographic description in the repository

*Publisher copyright*

(Article begins on next page)

This is an electronic version (author's or accepted version) of the paper:

Lamberti F., Sanna A., Demartini C., "A relation-based page rank algorithm for semantic web search engines," IEEE Transactions on Knowledge and Data Engineering, 21(1), pp. 123-136, 2009.

DOI: 10.1109/TKDE.2008.113

Link to IEEE Xplore®: [http://ieeexplore.ieee.org/xpls/abs\\_all.jsp?arnumber=4538222](http://ieeexplore.ieee.org/xpls/abs_all.jsp?arnumber=4538222)

© 2009 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works.

# A Relation-Based Page Rank Algorithm for Semantic Web Search Engines

Fabrizio Lamberti, *Member, IEEE*, Andrea Sanna, and Claudio Demartini, *Member, IEEE*

**Abstract**—With the tremendous growth of information available to end users through the Web, search engines come to play ever a more critical role. Nevertheless, because of their general purpose approach, it is always less uncommon that obtained result sets provide a burden of useless pages. Next generation Web architecture, represented by Semantic Web, provides the layered architecture possibly allowing to overcome this limitation. Several search engines have been proposed, which allow to increase information retrieval accuracy by exploiting a key content of Semantic Web resources, that is relations. However, in order to rank results, most of the existing solutions need to work on the whole annotated knowledge base. In this paper we propose a relation-based page rank algorithm to be used in conjunction with Semantic Web search engines that simply relies on information which could be extracted from user query and annotated resource. Relevance is measured as the probability that retrieved resource actually contains those relations whose existence was assumed by the user at the time of query definition.

**Index Terms**—Semantic Web, Knowledge retrieval, Search process, Query formulation



## 1 INTRODUCTION

IN the last years, with the massive growth of the Web we assisted to an explosion of information accessible to Internet users. Nevertheless, at the same time it has become ever more critical for end users to explore this huge repository and find needed resources by simply following the hyperlink network as foreseen by Tim Berners-Lee in 1999 [4]. Today, search engines constitute the most helpful tools for organizing information and extracting knowledge from the Web [9]. However, it is not uncommon that even the most renowned search engines return result sets including many pages which are definitely useless for the user [18]. This is mainly due to the fact that the very basic relevance criterions underlying their information retrieval strategies rely on the presence of query keywords within the returned pages. It is worth observing that statistical algorithms are applied to “tune” the result and, more important, approaches based on the concept of relevance feedback are used in order to maximize the satisfaction of user’s needs. Nevertheless, in some cases this does not suffice.

In order to show this oddly effect, let us see what happens when a user enters a query composed by the following keywords “hotel”, “Rome” and “historical centre” (or “hotel”, “Roma” and “centro storico”) in the Italian version of the well-known Google search engine. He would not be astonished probably by finding that the result set actually includes several hotels located in the historical centre of Rome, as expected. Another hotel located in a small town at some distance from Rome city centre is also included. However, two hotels located in

the historical centre of other main Italian cities are also displayed. Finally, three hotels named Roma are included among the ten most relevant results even if they have nothing to do with the selected city. Only four out the ten results presented to the user satisfy user needs (even if they seem to satisfy user query, based on the strategy adopted to process it). There is no doubt that the user would be able to easily decide which results are really of interest by looking for example at the two-line excerpt of the Web page presented in the displayed list or by quickly examining each page. Anyway, the presence of unwanted pages in the result set would force him to perform a post-processing on retrieved information to discard unneeded ones. Even though several automatic techniques have been recently proposed [32], result refinement remains a time waste and click expensive process, that is even more critical when the result set has to be processed by automatic software agents. Let us try to analyze more in detail the reason why “out-of-scope” pages are inserted in the result set. When the user entered the query “hotel”, “Rome” and “historical centre” he was assuming the existence of some relations among those terms, such as for example “hotel” *located in the* “historical centre” *of* “Rome”. However, when the query was sent to the search engine logic, these hidden details were lost. Search logic usually tries to recover this information by exploiting many text-matching techniques (such as number of occurrences and distance among terms). Nevertheless, traditional search engines do not have the necessary infrastructure for exploiting relation-based information that belongs to the semantic annotations for a Web page.

Semantic Web [5] will offer the way for solving this problem at the architecture level. In fact, in Semantic Web each page possesses semantic metadata which record additional details concerning Web page itself. Annotations are based on classes of concepts and relations among

• Fabrizio Lamberti, Andrea Sanna and Claudio Demartini are with the Dipartimento di Automatica ed Informatica, Politecnico di Torino, C.so Duca degli Abruzzi 24, 10129 Torino, Italy. E-mail: lamberti@polito.it, sanna@polito.it, demartini@polito.it.

<sup>1</sup> Discussion based on result set obtained on August 4, 2007.

them. The “vocabulary” for the annotation is usually expressed by means of an ontology that provides a common understanding of terms within a given domain.

In this paper, we will prove that relations among concepts embedded into semantic annotations can be effectively exploited to define a *ranking strategy* for Semantic Web search engines. This sort of ranking behaves at an inner level (that is, it exploits more precise information that can be made available within a Web page) and can be used in conjunction with other established ranking strategies to further improve the accuracy of query results. With respect to other ranking strategies for Semantic Web, our approach only relies on the knowledge of user query, Web pages to be ranked and underlying ontology. Thus, it allows to effectively manage the search space and to reduce the complexity associated with the ranking task.

The organization of this paper is as follows. In Section 2 we provide an overview of existing strategies for Semantic Web search. In Section 3, the basic idea laying behind the proposed approach is presented by resorting to practical examples, while in Section 4 a formal methodology for deriving the general rule is illustrated. In Section 5, details concerning implementation are provided. An analysis of algorithm complexity is given in Section 6, while experimental results are discussed in Section 7.

## 2 RELATED WORKS IN SEMANTIC WEB SEARCH

The aim of this paper is to show how to make use of relations in Semantic Web page annotations with the aim of generating an ordered result set, where pages which best fit user query are displayed first. The idea of exploiting ontology based annotations for information retrieval is not new [7][8][11][25]. Nevertheless, these first works did not focus on semantic relations, which are considered (and expected) to play a key role in Semantic Web [16][24]. In fact, it has been recently outlined that in order to fully benefit of semantic contents, a way for achieving relation based ranking has to be found [2][16][19][26].

One of the first attempts to enhance Semantic Web search engines with ranking capabilities is reported in [19]. The authors define a similarity score measuring the distance between the systematic descriptions of both query and retrieved resources. They first explode an initial set of relations (properties) by adding hidden relations which can be inferred from the query. Similarity is then computed as the ratio between relation instances linking concepts specified in user query and actual multiplicities of relation instances in semantic knowledge base. This method is applied on each property, individually, and requires to explore all the Semantic Web instances. Moreover, the user is requested to specify all the relations of interest. Thus, since it is predictable that the number of relations will largely exceed the number of concepts [1], its applicability in real contexts is severely compromised. A similar approach, aimed at measuring the relevance of a semantic association (that is a path traversing several concepts linked by semantic relations) is illustrated in [26]. The authors provide an interesting definition of relevance as the reciprocal of the ambiguity of the association

itself. However, this approach suffers from the same limitations of [19], since queries have to be specified by entering both concepts and relations, and ambiguity is measured over each relation instance.

Nevertheless, the idea of exploring the set of relations which are implicit in user’s mind (but which are not made explicit in defining the query) has been pursued in many works. In [18] ontology-based lexical relations like synonyms, antonyms and homonyms between keywords (but not concepts) have been used to “expand” query results. In this case, search is targeted to Web, rather than to Semantic Web. In [27], a similar approach has been integrated into artificial intelligence methodologies to address the problem of query answering. In [3], query logs are used to construct a user profile to be later used to improve the accuracy of Web search. Semantic Web search from the point of view of user’s intent has been addressed also in [15][28], where the authors present two methodologies for capturing the user’s information need by trying to formalize its mental model. They analyze keywords provided during query definition, automatically associate related concepts, and exploit the semantic knowledge base to automatically formulate formal queries.

A slightly different methodology has been exploited in SemRank [2]. Here, the basic idea is still to rank results based on how predictable a result might be for the user, but based on how much information is conveyed by a result, thereby giving a sense of how much information a user would gain by being informed about the existence of the result itself. To achieve their goal, the authors define two measures named “uniqueness” and “discrepancy” which allow to account for specificity or deviation of a particular result with respect to instances stored in the database. An additional added value of SemRank is that in the computation of the ranking, it is exploited a so called “modulative relevance model” that is capable of taking into account the particular context/purpose in/for which a query has been submitted (conventional or discovery search). Even if the authors do not provide any analysis of the computational cost of their approach, it is reasonable to infer that, since to rank a single page information related to annotations of all the remaining pages is needed, performance of the proposed solution would hardly scale for huge Semantic Web environments.

An approach also based on the context and partially solving the problems above is taken in [26]. Here, the context (defined in this case as a subset of concepts and relations of the whole Semantic Web environment) of interest to the user, rather than specific concepts or relations, can be specified together with the query using an ad-hoc language. The authors assign “universal” and “user-defined” weights to each semantic relation/association, taking into account the context as well as other parameters like specificity and path-length. These weights are combined into a global formula where multiplying constants are specified by the user (or by an expert) and are strictly query dependent. Thus, in order to get accurate results, an intensive manual configuration step has to be performed, and this step cannot be valid for answering heterogeneous queries. A solution capable of partially overcoming limi-

tations above is presented in [22]. Here, a strategy for clustering concepts based on query keywords provided by the user is proposed. A spread activation process is applied to navigate the whole relation set and discover related nodes which could be of interest. This process is only partially automated since it has to be guided by information on knowledge provided by a domain expert.

A totally different solution is represented by OntoLook [16]. The basic idea is that if a graph based representation of a Web page annotation can be provided, where concepts and relations (together with their multiplicities) are modelled as vertices and weighted edges, respectively, it becomes possible to define a series of cuts removing less relevant concepts from the graph. This allows for the generation of a so called candidate relation-keyword set (CRKS) to be submitted to annotated database, which can significantly reduce the presence of uninteresting pages in the result set. It is worth observing that the strategy behind OntoLook only allows to empirically identify relations among concepts which should be less relevant with respect to user query. This information is used to reformulate user query by including only a subset of all the possible relations among concepts, that is later used to retrieve Web pages from the annotated database. User is not requested to specify relations on interest during query definition. However, the effectiveness of the approach is strongly limited by the fact that there not exists any ranking strategy. Even if the authors claim that any of the existing page ranking algorithms can be used to order the obtained result set, it is worth remarking that this is not completely true. In fact, a ranking strategy like the PageRank [13][17] used by Google [6] is only one of the ranking algorithms used to organize results to be displayed to the user. Many other statistical and text-matching techniques are used together with PageRank. Of course, PageRank can be used in conjunction with [16] to exploit relevance feedback and post-process the result set. But, the use of the remaining techniques is not feasible since they cannot be reasonably applied into a concept-relation based framework where ontology is predominant on pure text. The authors themselves state that what is really needed is a relation-based page rank algorithm.

Our work moves from the considerations above and relies on the assumption that for providing effective ranking, search engine logic should only need to know the structure of the underlying ontology and of the Web page to be ranked in order to compute the corresponding relevance score. In this way, effective performance can be achieved in heterogeneous real frameworks. It is worth observing that the proposed approach could be easily seen as an extension of [16]. Moreover, it does not represent an alternative to any of the approaches above, but rather, they can be regarded as complementary to our solution (and vice versa). For instance, the availability of an ad-hoc language allowing the user to pre-process the graph and reduce the region of interest [26] could be integrated in our approach as a pre-processing step. Similarly, the availability of instruments for inferring concepts of interest starting from a pure keyword-based query [22] can be helpful to limit the amount of knowledge of the

underlying ontology requested to the user. Finally, the proposed technique is not intended to replace the ranking strategies of actual search engines. In fact, it relies on relevance information which are totally different from those exploited for example in algorithms like SemRank, Pagerank, and others. Rather, it should be understood as a pre-processing step to produce a semantic aware ordered result set to be later (or simultaneously) treated with existing (popular) techniques in order to come to an increased hit ratio in user query processing.

### 3 OVERVIEW OF THE RANKING STRATEGY

In this section, basic idea behind our ranking strategy is discussed. In order to introduce the readers to its formalism and let them foresee its applicability in real scenarios, the overall architecture of a prototypal search environment developed in our laboratory is presented first.

#### 3.1 Prototype of a Relation-based Search Engine

To evaluate the feasibility of the proposed approach, we first constructed a controlled Semantic Web environment. To do this, we selected the well-known *travel.owl* ontology [20] written in the OWL language [29] and we modified it by adding new relations in order to make it more suitable for demonstrating system functionality. We then created a knowledge base by either downloading or automatically generating a set of Web pages in the field of tourism and we embedded into them RDF [21] semantic annotations based on the ontology above. Finally, we designed the remaining modules of the architecture including a Web-page database, a crawler application, a knowledge database, an OWL parser (OwlDotNetApi), a query interface and the true search engine module embedding the proposed ranking logic (Fig. 1). The crawler application collects annotated Web pages from the Semantic Web (in this case represented by the controlled environment and its Web page collection) including RDF metadata and originating OWL ontology. RDF metadata are interpreted by the OWL parser and stored in the knowledge database. A graphics user interface allows for the definition of a query which is passed on to the relation-based search logic. Ordered result set generated by this latter module is finally presented to the user. Details of system workflow will be provided in the following sub-sections, starting with query definition process, since it was through the analysis its dynamics that we came to the identification of our ranking strategy.

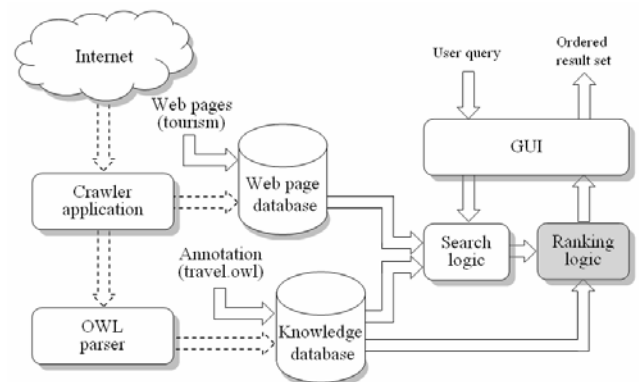


Fig. 1. Semantic Web infrastructure (prototype architecture).

### 3.2 Starting Point: the Query Definition Process

In a traditional search engine like Google [6], a query is specified by giving a set of keywords, possibly linked through logic operators and enriched with additional constraints (i.e. document type, language, etc.). On the other hand, semantic search engines are capable of exploiting concepts (and relations) hidden behind each keyword together with natural language interpretation techniques to further refine the result set. The core step that consists in identifying the mapping between keywords and concepts can be performed in a (semi-) automated way [15][22][23][28]. Otherwise, in order to avoid ambiguities, the user can be requested, during query definition, to specify the concept a keyword refers to [2][16][19][26]. Given the fact that the query interpretation step is out of the scope of this paper, the proposed methodology relies on the second approach. That is, like in [16], the user specifies a query by entering a keyword and selecting a concept from a pull-down menu containing ontology classes of the *travel.owl* ontology organized in a hierarchical fashion.

It is worth observing that the current implementation is not able to handle multiple ontologies describing the same domain. From the point of view of the search logic, this would require the integration of one of the existing techniques for mapping or merging/translating the heterogeneous ontologies [10] which would result in the definition of a set of mapping rules or in the creation of a novel (possibly extended) ontology, respectively. From the point of view of user interaction, having an extended ontology would increase the need for a pre-processing step enabling automatic identification of keyword-concept pairs. On the other hand, mapping rules would have to be only implemented in the search logic; the user could continue to use the same interface (possibly allowing to choose the ontology best suited for the query).

### 3.3 Introduction to Relation-based Ranking

Let us assume now that the user specifies the keyword “Rome” and he then selects from the pull-down menu one of the possible concepts such as *Destination*, or *City*. A second keyword “hotel” is then added, choosing *Accommodation* as the associated concept. In general, there is no way to state which was the relation in user’s mind between those two concepts (even if in this case it seems to be obvious). But, what it can be certainly said is that the user was assuming the existence of at least one relation between the two terms (and concepts, as well) or between these terms and the following ones (if this is the case). Now, let us consider a set of annotated pages containing keywords “Rome” and “hotel” and associated concepts *Destination* and *Accommodation*. A traditional search engine like Google would return both pages without considering the information provided by the semantic mark. On the other hand, a semantic search engine would take into account associations keyword-concept and would return a page only if both keywords (or synonyms, homonyms, etc.) are present within the page and they are related to associated concepts. Finally, a relations-based search

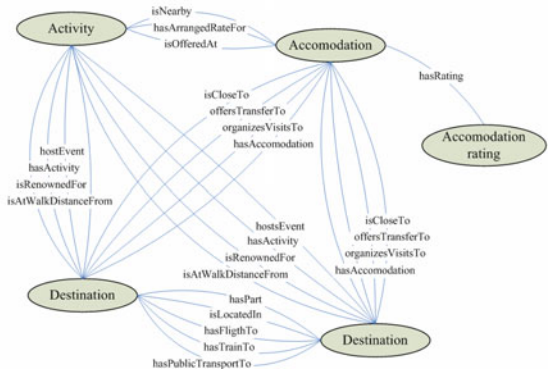


Fig. 2. A portion of the graph-based representation for *travel.owl* ontology (ontology graph).

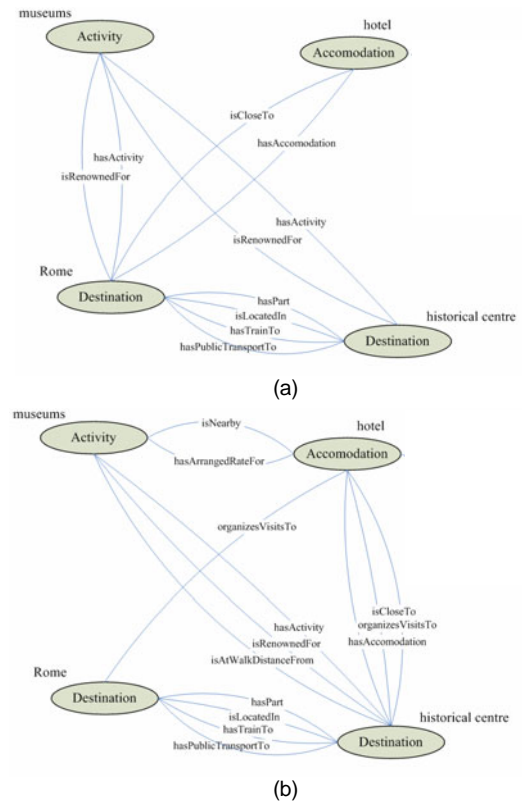


Fig. 3. Example of annotation graph for two Web pages. (a) Activities, accommodations and sightseeing places in Rome. (b) Hotel in the historical centre of Rome, close to museums.

engine like the one presented in [16] would go beyond pure “keyword isolated” search and would include these pages in the result set only if there exist enough relations linking considered concepts. However, pages included in the result set would have the same “weight”.

### 3.4 Basic Idea

Let us try to see if there is a way for presenting these pages in order of importance to the user. We continue with query definition and we assume now that the user enters the last two keywords of his query, like for example “museums” and “historical centre”, associated to concepts *Activity* and *Destination*, respectively. Let us assume also that, according to the ontology, these concepts are linked to both the previous concepts through a certain number of relations. There is again no way to infer either to which concept/s and by means of which relation/s the

TABLE 1  
DEFINITION OF SYMBOLS

Symbol	Definition
$G(C, R)$	Ontology graph
$C = \{c_1, c_2, \dots, c_n\}$	Set of concepts constituting the vertices of the ontology graph
$R = \{R_{ij} \mid i=1, \dots, n, j=1, \dots, n, j > i\}$	Set of relations constituting the edges of the ontology graph
$R_{ij} = \{r_{ij}^1, r_{ij}^2, \dots, r_{ij}^m, m < n\}$	Set of relations between concepts $c_i$ and $c_j$ in $G(C, R)$
$Q = \{(k_i, c_i)\}$	Query as a collection of pairs (keyword, concept)
$G_Q(C_Q, R_Q)$	Query sub-graph for query $Q$
$C_Q = \{c_i \mid (k_i, c_i) \in Q\}$	Set of concepts constituting the query sub-graph
$R_Q = \{\bar{R}_{ij} \mid 1 \leq i \leq n, 1 \leq j \leq n, j > i\}$	Set of relations constituting the query sub-graph
$\bar{R}_{ij} = \{\bar{r}_{ij} \mid c_i, c_j \in C_Q,  R_{ij}  \geq 1\}$	Set of relations between $c_i$ and $c_j$ in the query sub-graph
$\eta_{ij} =  \bar{R}_{ij} $	Number of relations between $c_i$ and $c_j$ in $G_Q(C_Q, R_Q)$
$A = (AC, AR)$	Graph-based annotation ( $AC$ and $AR$ are the sets of concepts and relations)
$AR_{ij} = \{r_{ij}^d \mid r_{ij}^d \in AR, 1 \leq d \leq m\}$	Set of relations between $c_i$ and $c_j$ in $AR$
$G_{Q,p}(C_{Q,p}, R_{Q,p})$	Page sub-graph for page $p$ given the query $Q$
$C_{Q,p} = \{c_i \mid c_i \in C_Q \cap AC\}$	Set of concept of page sub-graph $G_{Q,p}(C_{Q,p}, R_{Q,p})$
$R_{Q,p} = \{\bar{r}_{ij} \mid c_i, c_j \in C_{Q,p}\}$	Set of relations of page sub-graph $G_{Q,p}(C_{Q,p}, R_{Q,p})$
$\delta_{ij} =  AR_{ij} $	Number of relations $c_i$ and $c_j$ in $G_Q(C_Q, R_Q)$
$\tau_{ij} = P(\bar{r}_{ij}, p) = \delta_{ij} / \eta_{ij}$	Relation probability for $\bar{r}_{ij}$ in page $p$ given the query $Q$
$SF_{Q,p}(l)$	Set of spanning forests ( $l$ edges) for page $p$ and query $Q$
$SF_{Q,p}^f(l)$	$f$ -th spanning forest ( $l$ edges) for page $p$ and query $Q$
$\sigma_{Q,p}(l) =  SF_{Q,p}(l) $	Number of spanning forests ( $l$ edges) for page $p$ and query $Q$
$P(Q, p, l)$	Constrained relevance score for page $p$ , query $Q$ , relevance class $l$
$ps_{Q,p}$	Relevance score of page $p$ for a given query $Q$

newly added concepts are related to. However, we can certainly say that, since these are the last concepts, they should be related each other or to at least one of the previously entered concepts. In general, what we can always say is that *each concept specified within the query should have to be characterized by relations with at least another concept*. This consideration can be of great help when trying to define a way for providing a ranking among semantic annotated pages. In fact, the larger is the number of relations linking each concept with each other concept given the total number of relations among those concepts in the ontology, the higher is the probability that this page contains exactly those relations which are of interest to the user and, as a consequence, that this page is actually the most relevant with respect to user query. Thus, the idea is to define a “ranking criterion” based on *an estimate of the probability that keywords/concepts within an annotated page are linked one to the other in a way that is the same (or at least that is similar) to the one in user’s mind at the time of query definition*. As it will be shown in the following, this probability measure can be effectively computed by defining a graph-based description of the ontology (*ontology graph*), of the user query (*query sub-graph*), as well as of each annotated page containing queried concepts/keywords (both in terms of *annotation graph* and *page sub-graph*). In the following, ontology graph, query sub-graph, annotation graph and page sub-graph notions will be presented through the use of intuitive examples.

### 3.5 Graph-based Notation and Methodology

In the ontology and annotation graphs, concepts and relations are translated into graph nodes and edges, respec-

tively. A portion of the ontology graph for the *travel.owl* ontology is reported in Fig. 2, while two examples of annotation graphs built upon as many annotated Web pages are shown in Fig. 3. It is worth observing that, by taking into account the considerations in Section 3.4, a ranking for Web pages in Fig. 3 can be easily found: in fact, according to the query, the user was probably looking for a hotel located in the historical centre of Rome and (maybe) close to museums. However, even if this ranking can be proved intuitively by looking at the actual relations, a way for instructing the logic of the search engine is still needed. To do this, the notions of query sub-graph and page sub-graph have to be introduced.

In a query sub-graph, nodes are represented by concepts which have been specified within the query. Nodes/concepts are linked by a (weighted) edge only if there exists at least one relation between those concepts in the ontology. The weight is represented by the actual number of relations. Similarly, a page sub-graph is built based on the annotation associated to page itself.

The methodology we propose in this paper starts from page sub-graph computed over an annotated page and generates all the possible combinations of the edges belonging to the sub-graph itself non including cycles. Since there could exist pages in which there are concepts which do not show any relations with other concepts but that could still be of interest to the user, the methodology progressively reduces the number of edges in page sub-graph and computes the probability that each of the resulting sub-graphs obtained by a combination of the remaining edges is the one which matches user’s intention. Edge removal could lead to having concepts without



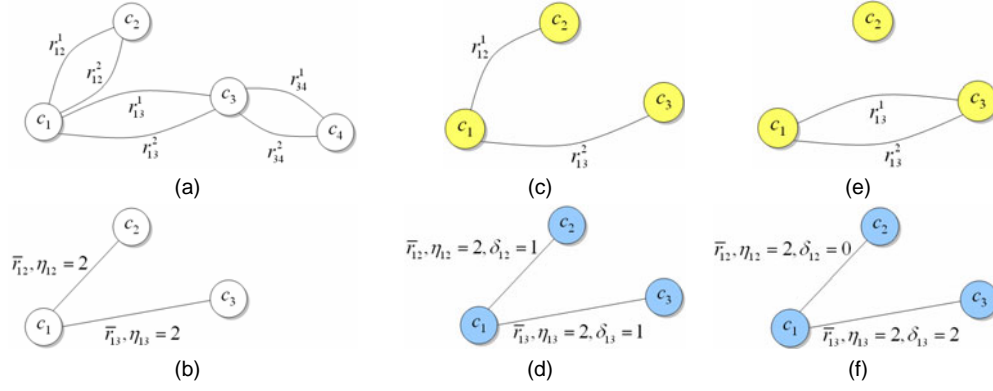


Fig. 4. (a) An ontology graph. (b) Query sub-graph obtained for given query specifying concepts  $c_1$ ,  $c_2$  and  $c_3$ . (c)-(d) A first example of page annotation  $p_1$  and related page sub-graph. (e)-(f) A second example of page annotation  $p_2$  and related page sub-graph.

any relation with other concepts. Thus, several relevance classes are defined, each characterized by a certain number of connected concepts. Within each class, pages are ordered depending on the probability measure above, and presented to the user.

## 4 RELATION-BASED RANKING FORMAL MODEL

In this section, a formal model for the proposed ranking strategy will be provided, by taking into account all the critical situations which could be envisioned

### 4.1 Graph-based Formalization

Starting from the ontology defined for a domain, a graph-based representation can be designed where OWL classes are mapped into graph vertices and OWL relation properties are mapped into graph edges. Thus, the existing relations between couples of concepts in the domain are depicted by means of connected vertices in the graph. We call it the *ontology graph*  $G$ . According to graph theory, the undirected graph  $G$  can be defined as  $G(C, R)$  where  $C = \{c_1, c_2, \dots, c_n\}$  is the set of concepts which can be identified in the ontology,  $|C| = n$  is the total number of concepts available,  $R = \{R_{ij} \mid i=1, \dots, n, j=1, \dots, n, j > i\}$  is the set of edges in the graph, and, more specifically,  $R_{ij} = \{r_{ij}^1, r_{ij}^2, \dots, r_{ij}^m, m < n\}$  is the set of edges between concepts  $i$  and  $j$ . An example of ontology graph (based on the formal notation summarized in Table 1) is illustrated in Fig. 4.a. Since queries are specified by the user by providing a collection of keywords and associated concepts, a single query can be formally expressed as  $Q = \{(k_i, c_i)\}$ .

Given a particular query containing a specific set of keywords related to a subset of ontology concepts, it is possible to construct a *query sub-graph*  $G_Q$ . Query sub-graph is an undirected weighted graph derived from  $G$  where vertices not belonging to  $C_Q$  are deleted. Moreover, in query sub-graph, vertices  $i$  and  $j$  are linked by an edge only if there exists at least one relation between the corresponding concepts in the ontology graph  $G$ . By referring to the same notation used for the ontology graph,  $G_Q$  can be expressed as  $G_Q(C_Q, R_Q)$ , where  $C_Q = \{c_i \mid (k_i, c_i) \in Q\} \subseteq C$  is the subset of concepts mentioned in the query,  $R_Q = \{\bar{r}_{ij} \mid 1 \leq i \leq n, 1 \leq j \leq n, j > i\}$  and  $\bar{R}_{ij} = \{\bar{r}_{ij} \mid c_i, c_j \in C_Q, |\bar{R}_{ij}| \geq 1\}$ . Each edge  $\bar{r}_{ij}$  in the query sub-

graph is assigned a weight  $\eta_{ij}$  that corresponds to the number of relations between concepts  $i$  and  $j$  in the ontology graph. Thus, it is  $\eta_{ij} = |\bar{R}_{ij}|$ . Query sub-graph that can be obtained from the ontology graph in Fig. 4.a for a query  $Q = \{(k_1, c_1), (k_2, c_2), (k_3, c_3)\}$  is shown in Fig. 4.b.

The aim of this paper is to demonstrate that, given an ontology graph  $G$  and a query sub-graph  $G_Q$ , it is possible to define a ranking strategy capable of assigning each page including queried concepts a relevance score based on the semantic relations available among concepts within the page itself (thus neglecting the contribution of the remaining Web pages). The proposed ranking strategy assumes that given a query  $Q$ , for each page  $p$  it is possible to build a *page sub-graph*  $G_{Q,p}$  using a methodology which is similar to the one used for  $G$  and  $G_Q$ , and exploiting information available in page annotation  $A$ . By expressing page annotation  $A$  as a graph we have  $A = (AC, AR)$ , where  $AC$  and  $AR$  are the sets of annotated concepts and relations, respectively. Page sub-graph  $G_{Q,p}$  contains only those concepts included both in  $C_Q$  and in page annotation  $AC$ . Concerning graph edges, all the edges  $\bar{r}_{ij}$  in  $R_Q$  are maintained, supposed that related concepts belong to  $G_{Q,p}$ . Weights  $\eta_{ij}$  specified for  $R_Q$  are inherited also by edges in  $G_{Q,p}$ . However, an additional weight  $\delta_{ij}$  is associated to each edge to take into account the number of relations actually linking concepts  $i$  and  $j$  in the selected page (on the basis of the set of annotated relations,  $AR$ ). According to the notation above, page sub-graph for page  $p$  can be defined as  $G_{Q,p}(C_{Q,p}, R_{Q,p})$  where  $C_{Q,p} = \{c_i \mid c_i \in C_Q \cap AC\}$  and  $R_{Q,p} = \{\bar{r}_{ij} \mid c_i, c_j \in G_{Q,p}\}$ . We have also  $\delta_{ij} = |AR_{ij}|$ , where  $AR_{ij} = \{r_{ij}^d \mid r_{ij}^d \in AR, 1 \leq d \leq m\}$ .

### 4.2 Relevance and Semantic Relations

Let us pass now at considering how to apply the methodology above for the computation of a page relevance score. We start again by analyzing (now from a formal point of view) the steps followed by a user during the process of query definition. Let us imagine that user is interested in pages containing three generic keywords  $k_1$ ,  $k_2$  and  $k_3$  (associated to as many generic concepts  $c_1$ ,  $c_2$  and  $c_3$ ). User begins query definition by specifying a pair including a keyword and its related concept. Let us assume that he starts with  $k_1$  and  $c_1$ . It is reasonable to assume that, after specifying keyword  $k_1$ , the user inserts



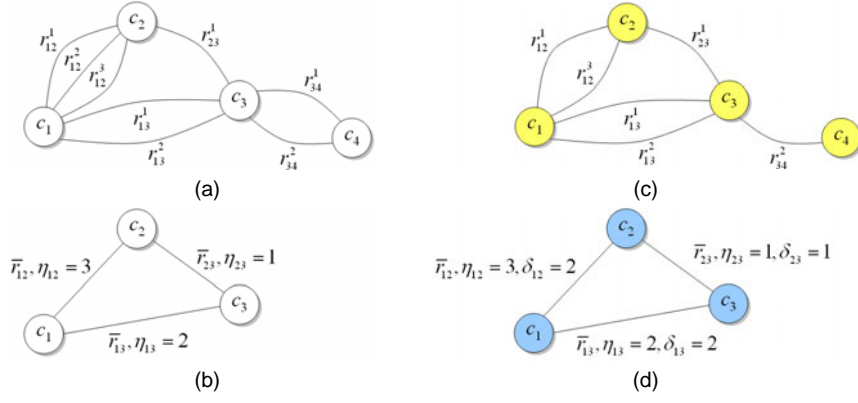


Fig. 5. (a) An ontology graph. (b) Query sub-graph. (c) An example of annotated page. (d) Page sub-graph built upon given ontology/query.

a second keyword (for example  $k_2$ , together with concept  $c_2$ ) expecting either to find pages where  $k_1$  and  $k_2$  (that is,  $c_1$  and  $c_2$ ) are related in some way or to find pages where  $k_1$  is linked to some other keywords/concepts that will be specified later. In a similar way, when he specifies  $k_3$  and  $c_3$  he would be expecting to further adjust the result set in order to find pages showing also relations between  $k_3$  and  $k_1$  (not  $k_2$  since in the ontology there is no relation linking  $c_3$  with  $c_2$ ). Let us consider a very trivial example assuming that there exists only two pages  $p_1$  and  $p_2$  containing all the keywords (and associated concepts) specified by the user. This represents the (initial) result set for the given query. We want to rank those pages in order to present to the user first the page that best fits his query. Semantic annotations and page sub-graphs for these pages are illustrated in Fig. 4.c-f. In the first page, both  $c_2$  and  $c_3$  are linked to  $c_1$  through a single relation (Fig. 4.c), while in the second page there exists two relations linking  $c_3$  to  $c_1$ . However,  $c_2$  is not linked in any way to  $c_1$  (Fig. 4.f) Since we cannot assume which could be the concepts or the relations more important with respect to user query, we can provide a significant measure of page relevance by computing the probability that a page is the one of interest to the user (that is, its relevance) by calculating the probability that  $c_2$  is linked to  $c_1$  and  $c_3$  is linked to  $c_1$  through the relations in user's mind (either  $r_{12}^1$  or  $r_{12}^2$  and  $r_{13}^1$  or  $r_{13}^2$ , respectively). Let us compute  $P(\bar{r}_{ij}, Q, p)$ , that is the probability of finding in a particular page  $p$  a relation  $\bar{r}_{ij}$  between concepts  $i$  and  $j$  which could be the one of interest to the user (because of query  $Q$ ). According to the probability theory, this can be defined as  $P(\bar{r}_{ij}, p) = \delta_{ij} / \eta_{ij} = \tau_{ij}$  (note that it does not depend on  $Q$ ). We call it the *relation probability*. Thus, for the first page we have  $P(\bar{r}_{12}, p_1) = \delta_{12} / \eta_{12} = \tau_{12} = 1/2$  and  $P(\bar{r}_{13}, p_1) = \delta_{13} / \eta_{13} = \tau_{13} = 1/2$ . For the second page we have  $P(\bar{r}_{12}, p_2) = \delta_{12} / \eta_{12} = \tau_{12} = 0$  and  $P(\bar{r}_{13}, p_2) = \delta_{13} / \eta_{13} = \tau_{13} = 1$ . Based on the considerations above, we can compute the joint probability  $P(Q, p) = P((\bar{r}_{12}, p) \cap (\bar{r}_{13}, p))$ . The dependency on  $Q$  is due to the fact that only concepts given in  $Q$  are taken into account. Since the events  $(\bar{r}_{12}, p)$  and  $(\bar{r}_{13}, p)$  are not correlated,  $P(Q, p)$  can be rewritten as  $P(Q, p) = P(\bar{r}_{12}, p) \cdot P(\bar{r}_{13}, p)$ . Thus, for the specific example being considered, it is  $P(Q, p_1) = 1/4$  and  $P(Q, p_2) = 0$ , respectively for the first and second page. This allows to place the first page before the second one in the ordered

result set. However, to preserve the behaviour of common search strategies, a way for assigning a score different than zero to pages in which there exists concepts not related to other concepts will have to be identified.

Another critical situation is illustrated in Fig. 5. In this case, the user specifies a query composed by concepts  $c_1$ ,  $c_2$ , and  $c_3$  over a novel ontology. Based on the considerations above, a measure of page relevance can be computed by estimating, for each concept, the probability of having a relation between that concept and another concept and that such relation is exactly the one in user's mind. However, it can be demonstrated that this probability can be expressed also in different terms, capable of taking into account situations in which a particular concept can be related to more than one concept (that is the case of the specific example being considered as well as of common situations in any concrete search scenario). Specifically, the probability that each concept is related to other concepts is given by the probability of having  $c_1$  linked to  $c_2$  and  $c_2$  linked to  $c_3$ , or  $c_1$  linked to  $c_2$  and  $c_1$  linked to  $c_3$ , or  $c_2$  linked to  $c_3$  and  $c_1$  linked to  $c_3$ . The situations above can be modelled again by using graph theory. In fact, having each concept related to at least another concept in the query is equivalent to considering all the possible spanning forests (a collection of spanning trees, one for each connected component in the graph) for page sub-graph  $G_{Q,p}$  given the query  $Q$ . In Fig. 6, all the possible spanning forests (trees, in this case) of the page sub-graph in Fig. 5.d are shown. We call  $SF_{Q,p}^f$  the  $f$ -th *page spanning forest* computed over  $G_{Q,p}$ . We define  $P(SF_{Q,p}^f)$  as the probability that  $SF_{Q,p}^f$  is the spanning forest of interest to the user. By simplifying the notation and replacing  $\bar{r}_{ij}, p$  with  $\bar{r}_{ij}^p$ , the probability for page  $p$  can be computed as:

$$P(Q, p) = P\left(\left((\bar{r}_{12}^p \cap \bar{r}_{23}^p) \cap SF_{Q,p}^1\right) \cup \left((\bar{r}_{12}^p \cap \bar{r}_{13}^p) \cap SF_{Q,p}^2\right) \cup \left((\bar{r}_{23}^p \cap \bar{r}_{13}^p) \cap SF_{Q,p}^3\right)\right). \quad (1)$$

Since the events are not correlated, it is also:

$$\begin{aligned} P(Q, p) &= P(\bar{r}_{12}^p \cap \bar{r}_{23}^p) \cdot P(SF_{Q,p}^1) + P(\bar{r}_{12}^p \cap \bar{r}_{13}^p) \cdot P(SF_{Q,p}^2) \\ &\quad + P(\bar{r}_{23}^p \cap \bar{r}_{13}^p) \cdot P(SF_{Q,p}^3) \\ &= P(\bar{r}_{12}^p) \cdot P(\bar{r}_{23}^p) \cdot P(SF_{Q,p}^1) + P(\bar{r}_{12}^p) \cdot P(\bar{r}_{13}^p) \cdot P(SF_{Q,p}^2) \\ &\quad + P(\bar{r}_{23}^p) \cdot P(\bar{r}_{13}^p) \cdot P(SF_{Q,p}^3) \end{aligned} \quad (2)$$

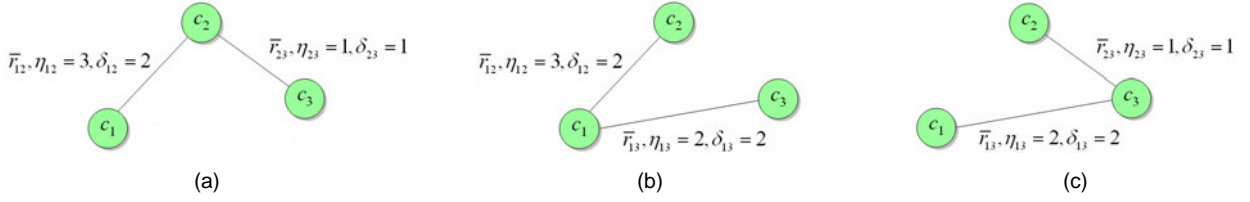


Fig 6. All the possible spanning forests (trees) which could be obtained from  $G_{Q,p}$  in Fig. 5.d.

where  $P(\bar{r}_{ij,p})$  can be replaced with  $\tau_{ij} = \delta_{ij} / \eta_{ij}$ .

Since the probability for a single page spanning forest to be the one of interest to the user is the same with respect to the remaining ones, if we define  $\sigma_{Q,p}$  as the number of spanning forests for  $G_{Q,p}$ , we have  $P(SF_{Q,p}^1) = P(SF_{Q,p}^2) = P(SF_{Q,p}^3) = 1/\sigma_{Q,p}$ . Thus, expression for  $P(Q,p)$  can be rewritten again as:

$$P(Q,p) = \frac{P(\bar{r}_{12}^p) \cdot P(\bar{r}_{23}^p) + P(\bar{r}_{12}^p) \cdot P(\bar{r}_{13}^p) + P(\bar{r}_{23}^p) \cdot P(\bar{r}_{13}^p)}{\sigma_{Q,p}} \quad (3)$$

and according to the definition of relation probability it is:

$$P(Q,p) = [\tau_{12} \cdot \tau_{23} + \tau_{12} \cdot \tau_{13} + \tau_{23} \cdot \tau_{13}] / \sigma_{Q,p} \quad (4)$$

Given the ontology and the query selected for the considered example, expression (4) can be used to compute a relevance score for each page in the result set and to provide a ranking within the result set itself. As expected, (4) works well also for the example in Fig. 4, where  $\sigma_{Q,p} = 1$  (since the page sub-graph already constitutes the only spanning forest). Nevertheless,  $P(Q,p)$  can still assume a value equal to zero for all those pages in which there exists concepts which do not show any relation with other concepts, but still present, as a keyword, in the annotated page. In the following, we will analyze this issue in details and we will show how to extend the methodology above in order to come to a general rule ranking all the pages in the (initial) result set.

We consider again an example represented by two pages (depicted in Fig. 7, and based on the same ontology in Fig. 5.a), where concept  $c_4$  (in the first page) and concept  $c_2$  (in the second page) do not show any relations with the remaining concepts. If we compute  $P(Q,p_1)$  and  $P(Q,p_2)$  using (4) (which is still valid since page annotation refers to the same ontology) we get a relevance score equal to zero. Based on the definition of relevance score provided above, in order to find a score different than zero allowing each page to be ranked with respect to other pages, we have to *relax the condition of having each concept related to each other concept*. Since, for definition, in a spanning forest does not exist any cycles, removing one edge means removing a link between a couple of concepts. That is, edges from all the page spanning forests have to be progressively removed, thus obtaining *constrained page spanning forests* composed by a decreasing number of edges (and equivalently, of connected concepts). We maintain the term “spanning” in order to recall that each constrained page spanning forest originates from a true spanning forest in which, for all the connected components of the graph, all the vertices are linked by exactly one edge. However, we introduce the term “constrained” to recall that there exists a constraint on the

number of edges of the forest allowing for the existence of not connected vertices in the graph. Since there is no way to infer which was the link between concepts more relevant to the user at the time of defining the query, constrained page spanning forests characterized by the same number of edges can be considered as comparable in terms of relevance with respect to user query. All the constrained page spanning forests composed by the same number of edges represent a possible (even if less relevant) answer to user query. Based on the number of constrained page spanning forests that can be generated from the page sub-graph for a given number of edges, the probability of that page can be calculated as the sum of the probabilities computed for each constrained page spanning forest of a given length divided by the total number of constrained page spanning forests of that length that can be originated by the page sub-graph. In the following, this latter consideration together with empirical results presented in this section will be exploited in order to provide a general rule for relation-based ranking of semantic annotated Web pages.

### 4.3 Page Relevance Score and Ranking

Let us consider an ontology graph  $G$  and a query sub-graph  $G_Q$ . Let us consider a page  $p$  and let us derive its page sub-graph  $G_{Q,p}$ . We now define  $SF_{Q,p}^f(l)$  that is the set including all the constrained spanning forests for a given number of edges  $l$  ( $1 \leq l < C_{Q,p}$ ). The cardinality of this set is  $\sigma_{Q,p}(l) = |SF_{Q,p}^f(l)|$ . Finally, let us define  $SF_{Q,p}^f(l)$  as the  $f$ -th spanning forest originated from the page sub-graph for the given query  $Q$  and page  $p$  and a specific number of edges  $l$ . When  $l$  is equal to the maximum length of a spanning forest of the page sub-graph, this correspond to a page spanning forest. Otherwise, it corresponds to a constrained page spanning forest. To simplify the explanation, we will sometimes refer to both these forests as page forests (except when this can cause ambiguities). The probability that a page forest  $SF_{Q,p}^f(l)$  is the one of interest to the user can be written as  $P(SF_{Q,p}^f(l))$ . Based on the considerations above, it is possible to define a constrained relevance score for page  $p$  as:

$$P(Q,p,l) = P \left( \bigcup_{f=1}^{|SF_{Q,p}^f(l)|} \left( \bigcap \{ \bar{r}_{ij,p} \mid \bar{r}_{ij,p} \in SF_{Q,p}^f(l) \} \cap SF_{Q,p}^f(l) \right) \right) \quad (5)$$

$$= \sum_{f=1}^{|SF_{Q,p}^f(l)|} \prod_{\bar{r}_{ij,p} \in SF_{Q,p}^f(l)} P(\bar{r}_{ij,p}) \cdot P(SF_{Q,p}^f(l))$$

where  $P(SF_{Q,p}^f(l)) = 1/\sigma_{Q,p}(l)$ . We call it a *constrained page relevance score* since its value depends on the value of  $l$ .

By iteratively considering all the constrained

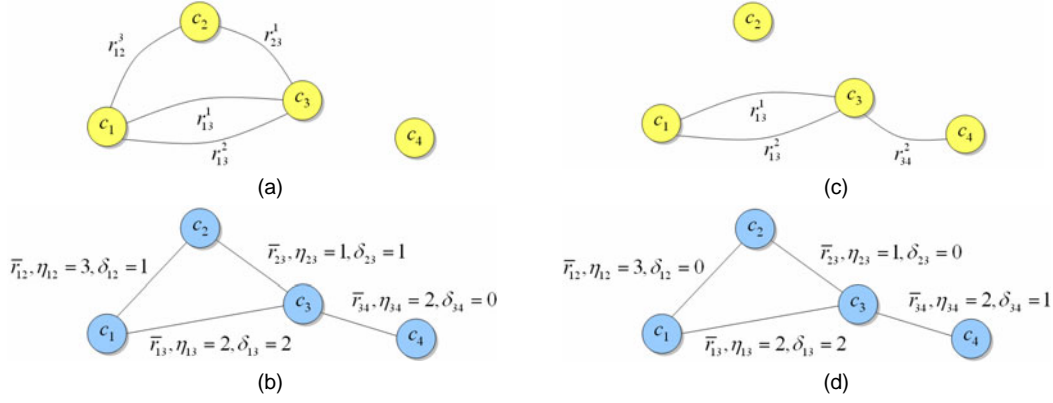


Fig. 7. (a) An annotated page  $p_1$  where concept  $c_4$  is not linked to any other concepts. (b) Page sub-graph for a query  $Q$  specifying  $c_1, c_2, c_3$  and  $c_4$ . (c) Annotation of a second page  $p_2$  where  $c_2$  is not linked to any other concepts. (d) Page sub-graph for the same query.

spanning forests characterized by the same length, we are progressively relaxing the constraint of having all the concepts related in some way to other concepts within the page. As soon as a value different than zero is obtained for  $P(Q, p, l)$ , we assume that this corresponds to a “final” relevance score for that page.

However, since  $P(Q, p, l)$  is computed as a probability, we have  $0 \leq P(Q, p, l) \leq 1$ . Thus,  $P(Q, p, l)$  cannot be directly used to compare one page in the result set with the remaining ones. Nevertheless, we can exploit the information on  $l$  to create several relevance classes in a straightforward way. In fact, reducing the value of  $l$ , as soon as we find a value different than zero for  $P(Q, p, l)$ , we compute the *page relevance score* (or *page score*) as:

$$ps_{Q,p} = P(Q, p, \max(l)) + \max(l) \mid P(Q, p, l) \neq 0 \quad (6)$$

In this way, each relevance class contains pages with a score in the range  $[l, l+1]$ , and pages within the same class are directly comparable and the (final) result set can be ordered by decreasing values of page score.

## 5 IMPLEMENTATION OF THE RANKING ALGORITHM

### 5.1 Overall Procedure

We now assemble the various steps illustrated in previous sections to present the overall ranking methodology (whose workflow is depicted in Fig. 8). User starts defining query keywords and concepts. Search engine logic accesses the Web page database, constructs the initial result set including all those pages which contain queried keywords and concepts, and computes the query sub-graph. Then, for each page in the result set, page sub-graph is computed. Starting from each sub-graph, all page spanning forests (both constrained and unconstrained) are generated and used to compute page score based on (6). Web pages are associated to relevance classes and the final (ordered) result set is constructed.

### 5.2 Spanning Forest Generation Algorithm

According to (6), calculating the relevance score for a single page requires to consider all the page forests and, for each forest, to compute the constrained page relevance score. This requires to find an efficient way for both enumerating all the page forests for a given query and com-

puting page probability.

Two strategies are feasible in order to approach the problem above. The first strategy could be to consider all the possible page spanning forests (page forests including a number of edges equal to the number of nodes minus one) of the page sub-graph and progressively remove their edges generating constrained page spanning forests by taking care to avoid duplicate configurations. In the worst case, all the edges have to be recursively removed until page forests with a single edge are generated. It is worth observing that from a computational point of view, avoiding the production of duplicate configurations is an extremely expensive task. A second strategy could start considering all the page forests of length one and generate all the possible page forests of increasing length by recursively adding a new edge until a page spanning forest is obtained. With respect to the previous approach, this second method shows several advantages. First, by properly selecting the edge to be added in the recursive process it is possible to implicitly obtain a set of page forests without duplicates. Moreover, the iterative approach allows to exploit results achieved in previous steps in order to speedup the time requested for computation. In fact, the probability associated with a particular forest made up of a given number of edges can be obtained by simply taking into account the contribution of the newly added edge. That is, computing the probability of a forest composed by  $n$  edges simply requires to multiply the probability obtained for page forest with  $n-1$  edges by the relation probability associated to the additional edge.

Thus, in this work we chose the second approach. Unfortunately, even if many algorithms have been proposed in the literature for addressing the task of finding all the spanning forests (or trees) in a graph [14][30], none of them is capable of taking into account forests with a variable number of edges derived from originating spanning forests. Thus, an ad-hoc algorithm has been designed (whose pseudo code is reported in Fig. 9). A detailed analysis of its complexity is provided in Section 6. It is worth observing that the incremental approach adopted in this algorithm shows an additional benefit with respect to the decremental one. In fact, it becomes possible to impose an upper bound to the growth of page forests in terms of number of edges. Since a larger number of edges

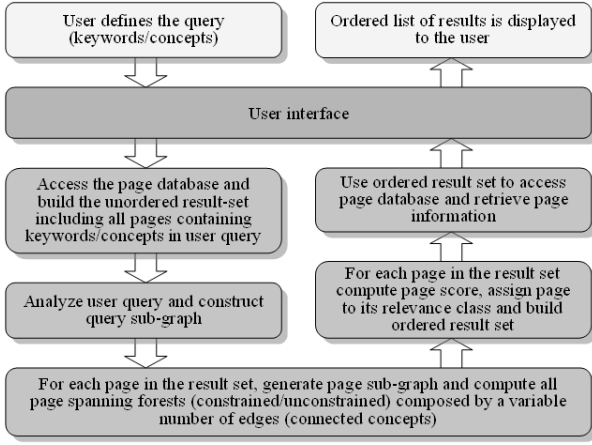


Fig. 8. Workflow, from query definition to presentation of results.

means a higher accuracy in the estimation of page relevance accompanied by a larger computational cost, the possibility of introducing a threshold to the widest page forest to be considered could allow to achieve a trade-off between ranking precision and complexity.

## 6 AN ESTIMATE OF ALGORITHM COMPLEXITY

### 6.1 Overview of the evaluation method

According to the ranking method presented above, the relevance score for a particular page is given by the first non-zero constrained relevance score obtained by varying the value of  $l$ . The number of times that length has to be varied depends on the annotation of the considered page. Moreover, the number of page forests for a given length depends on the topology of the page sub-graph. Thus, in order to provide an estimate of algorithm complexity, we have to consider the worst case, represented by a page based on an ontology with a complete graph and whose annotation includes only one relation. In this case, all the possible lengths have to be considered, and the maximum number of spanning forests has to be taken into account.

### 6.2 Worst Case Analysis

Let us consider an ontology characterized by a complete graph  $G$ . Let us assume that a query sub-graph  $G_Q$  has been defined and let us consider a page sub-graph  $G_{Q,p}$  where  $N_v = |C_{Q,p}|$  is the number of vertices and  $N_e = |R_{Q,p}|$  is the number of edges. For definition, given the number of edges of a complete page sub-graph, the number of  $l$ -length spanning forests  $\sigma_{Q,p}(l)$  is given by the number of sub-graphs composed by  $l$  edges not including cycles. In other words,  $\sigma_{Q,p}(l)$  can be computed as the overall number of possible combinations (without repetitions) of the  $N_e$  edges minus the number of sub-graphs composed by  $l$  edges and including a cycle. The number of  $l$ -combinations of  $N_e$  edges can be expressed as:

$$\kappa(l) = \text{bin}(N_e, l) \quad (7)$$

while the number  $\chi_\lambda(l)$  of  $l$ -length sub-graphs including an  $\lambda$ -cycle (with  $3 \leq \lambda \leq l$ ) and having  $l - \lambda$  additional “floating” (free) edges can be expressed as in (8). This expression takes into account the number  $\text{bin}(N_v, \lambda)$  of combinations of  $\lambda$  vertices over the  $N_v$  vertices of the

```

Label the edges in  $G_{Q,p}$  with an index ranging from 1 to  $R_{Q,p}$ 
Define variables  $e$  and  $a$  to index graph edges
Set  $\eta_e = \eta_{ij}$  //number of relations linking concepts
//  $i$  and  $j$  in the ontology graph (edge  $e \in R_{Q,p}$ )
Set  $\delta_e = \delta_{ij}$  //number of relations linking concepts
//  $i$  and  $j$  in the page sub-graph (edge  $e \in R_{Q,p}$ )
Set  $\tau_e = \eta_{ij} / \delta_{ij}$  //relation probability for edge  $e$ 
Mark all the edges in  $G_{Q,p}$  as not visited
Allocate weight vector  $W$  of size  $|C_{Q,p}| - 1$ 
//  $W[l]$  stores the accumulated constrained
// probabilities for page forests of length  $l$ 
Allocate vector  $\Sigma$  of size  $|C_{Q,p}| - 1$ 
//  $\Sigma[l]$  stores the number of page forests
// for a given length  $l$ 
Initialize  $W$  and  $\Sigma$  to zero

for  $e = 1, e \leq |R_{Q,p}|, e = e + 1$ 
    mark edge  $e$  as visited
    visit( $e, e, 1, \tau_e$ )
     $W[l] = W[l] + \tau_e$ 
     $\Sigma[l] = \Sigma[l] + 1$ 

function visit( $o, e, l, s$ )
     $a = e + 1$ 
    while  $a \leq |R_{Q,p}|$  and  $l \leq |C_{Q,p}| - 1$ 
        if  $a$  is not visited and  $a$  is safe
            // (does not introduce cycles, checked through DFS)
            mark edge  $a$  as visited
            visit( $o, a, l + 1, s \times \tau_a$ )
             $W[l + 1] = W[l + 1] + s$ 
             $\Sigma[l + 1] = \Sigma[l + 1] + 1$ 
            set edge  $a$  as not visited
        else
             $a = a + 1$ 

```

Fig. 9. Pseudo code of the algorithm for generating all the page spanning forests of variable length (incremental approach).

graphs as well as the possible  $(\lambda - 1)!/2$  configurations for each combination. For each combination and for each configuration, the presence of  $l - \lambda$  free edges which can be combined in  $\text{bin}(N_e - \lambda, l - \lambda)$  ways is considered. Since the presence of free edges could lead to the generation of cycles longer than  $\lambda$ , a correction factor is applied. Another correction factor is applied to deal with duplicate configurations possibly resulting into shorter cycles. Because of the novelty of this formulation, proofs for (8) will be provided in a specific paper<sup>2</sup>.

In conclusion, expression for  $\sigma_{Q,p}(l)$  can be written as:

$$\sigma_{Q,p}(l) = \kappa(l) - \sum_{\lambda=3}^l \chi_\lambda(l) \quad (9)$$

Values for  $\sigma_{Q,p}(l)$  computed using (9) over complete page sub-graphs including a number of concepts  $N_v$  ranging from two to seven are reported in Table 2. It is worth observing that, according to the methodology presented in previous section, computing the constrained relevance score requires  $\sigma_{Q,p}(l)$  accumulations plus one division for  $l=1$  and  $\sigma_{Q,p}(l)$  multiplications,  $\sigma_{Q,p}(l)$  accumulations and one division for  $1 < l \leq N_v - 1$ . In the worst case, ranking one page requires  $\sum_{l=1}^{N_v-1} \sigma_{Q,p}(l)$  accumulations,  $1 + \sum_{l=2}^{N_v-1} \sigma_{Q,p}(l)$  multiplications and  $N_v - 1$  divisions.

## 7 EXPERIMENTAL RESULTS AND COMPARISONS

In this section, the applicability of our technique into real

<sup>2</sup> An application to a concrete example is shown in the Appendix which is available on-line: <http://gohan.polito.it:8080/tkde/appendix.pdf>



TABLE 2  
NUMBER OF PAGE FORESTS FOR COMPLETE PAGE SUB-  
GRAPHS WITH INCREASING NUMBER OF NODES

	$N_V = 2$	$N_V = 3$	$N_V = 4$	$N_V = 5$	$N_V = 6$	$N_V = 7$
$l = 1$	1	3	6	10	15	21
$l = 2$		3	15	45	105	210
$l = 3$			16	110	435	1295
$l = 4$				125	1080	5250
$l = 5$					1296	13377
$l = 6$						16807

scenarios will be analyzed by conducting two types of evaluations aimed at measuring the performance in terms of both time complexity and accuracy. Time complexity will be compared with that of [16], since our technique could be easily seen as an extension of it. Nevertheless, since the methodology in [16] is not targeted at ranking the result set, the accuracy of results will be compared with that of a traditional search engine like Google.

### 7.1 Time Complexity

The computation of fair results concerning time complexity requires a sufficiently large repository with a significant number of annotated pages. Because of the difficulty of integrating the proposed technique within today's search engines like Google, in which a native semantic layer is actually missing, we chose to estimate the computation time over a synthetic Semantic Web environment. The positive effect of this choice is twofold. On one hand, it is possible to work on as many pages as needed, thus effectively simulating the next generation Semantic Web repositories. On the other hand, by statistically annotating Web pages, we do not incur in the risk of biasing the result. In order to compare our measures with those of [16], we worked with the same ontology (*travel.owl*) and we selected the same query (in the query, illustrated in Fig. 10, specific keywords and concepts defined in [16] have been replaced with numeric indexes). We automatically generated a Web page database with one million pages, each page containing all the keywords specified in the query. For each page we constructed a semantic annotation based on the concepts defined in the selected ontology, randomly associating to each keyword one of the concepts in the ontology. We adjusted the statistical parameters so as to obtain a set of approximately one hundred thousand pages (precisely, 96'843 pages) including at least one of the keywords associated exactly to the concept specified in the query. Finally, we added semantic relations between concepts by uniformly distributing them across pages. In this way, each pair of concepts was linked by a variable subset of the relations associated to that pair in the ontology (each page containing approxi-

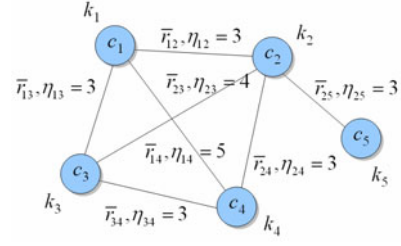


Fig. 10. Ontology from [16] used for measuring time complexity.

mately ten relations). Distribution of concepts and relations in the Web page database is summarized in Table 3. For each concept  $c_i | c_i \in C_Q$ , Table 3 reports the number of pages containing exactly the association keyword/concept  $(k_i, c_i)$  defined in the query. Moreover, for each couple of concepts  $c_i, c_j | c_i, c_j \in C_Q$ , it reports the number of pages in which both the concepts are associated to the keywords specified in user's query. Finally, for each couple of concepts  $c_i, c_j$ , Table 3 gives the number of pages which contain at least one of the relations  $\bar{r}_{ij}^d | \bar{r}_{ij}^d \in \bar{R}_{ij}, 1 \leq d \leq \eta_{ij}$  defined in the ontology, as well as the exact number of pages existing in the Web page database for each of the possible relations. Both the approach in [16] and the methodology presented in this paper have been applied onto the experimental environment above using an Intel Core 2 6400 CPU @ 2.13 GHz with 2GB of RAM. Results are shown in Table 4.

As illustrated in Section 2, in [16] a query graph is constructed over the ontology starting from concepts and keywords passed by the user. Then, the query graph is progressively reduced, thus obtaining several query sub-graphs (Table 4, column three). By taking for each edge in the query sub-graph one of the possible relations associated to that edge, several property-keyword pairs are generated (column four). The collection of these pairs constitutes a candidate relation-keyword set (CRKSs), that is submitted to the knowledge database for retrieving the page result set. With respect to [16], results related to CRKS generation present a speed-up due to the newer hardware used for the experiments and to an optimized procedure for combining contributing edges. However, overall delay (column six) is worse than in [16], since the time for submitting CRKSs to the database and for intersecting the results is also taken into account (column five).

Results obtained using our methodology are tabulated in columns from seven to nine. In particular, in column seven, the average number of page spanning forests for increasing length given the number of query concepts/keywords is reported. Moreover, column eight gives the time requested for extracting from the database the annotation of pages to be processed and for generating

$$\chi_\lambda(l) = \begin{cases} \binom{N_V}{\lambda} \cdot \frac{(\lambda-1)!}{2} & l = \lambda \\ \binom{N_V}{\lambda} \cdot \frac{(\lambda-1)!}{2} \cdot \left[ \binom{N_E - \lambda}{l - \lambda} - \left[ \frac{\lambda \cdot (\lambda-1)}{2} - \lambda \right] \cdot \binom{N_E - \lambda - 1}{l - \lambda - 1} \right] - \sum_{i=2}^{l-\lambda} \binom{N_V}{\lambda} \cdot \frac{(\lambda-1)!}{2} \cdot \frac{\lambda}{i} \cdot (N_V - \lambda) \cdot \binom{N_E - \lambda - i}{l - \lambda - i} & l > \lambda \geq 3 \\ - \lfloor (l - \lambda)/3 \rfloor \cdot \binom{N_V}{\lambda} \cdot \left( \frac{\lambda \cdot (N_V - \lambda)}{1 + \lfloor (l - \lambda)/3 \rfloor} + \frac{(N_V - \lambda)}{\lfloor (l - \lambda)/3 \rfloor} \right) \cdot \binom{N_E - \lambda - \lfloor (l - \lambda)/3 \rfloor \cdot \lambda}{l - \lambda - \lfloor (l - \lambda)/3 \rfloor \cdot \lambda} & \end{cases} \quad (8)$$

the associated page sub-graphs, together with the time requested for running the proposed algorithms and getting page relevance scores. It is worth observing that even if a final ordering is needed to sort the results, this delay has not been considered in Table 4 in order to let the user evaluate the time complexity of comparable result sets.

Experimental results show that the methodology in [16] is severely affected by the costs associated with the submission of CRKs to the database and with the intersection of results, while our spanning forest based approach allows to effectively manage the search space and to reduce the time complexity associated with the search task. The additional advantage of the proposed approach is that it incorporates the computation of a probability measure that can be effectively used to produce an ordered result set. Even if promising results over one million pages demonstrate the feasibility of the approach, at the same time they anticipate the need for further research activities aimed at ensuring scalability with the next generation Semantic Web repositories. In Section 5.2 we propose a practicable approach based on a threshold over the computed spanning forests' width. Nevertheless, we expect to further investigate this issue by analyzing the effect of the adoption of parallel and distributed computing paradigms, as well as of solutions for storing pre-computed (and updated) digests of page sub-graphs.

## 7.2 Accuracy

The accuracy of the proposed technique has been evaluated against the result set generated by running the query "hotel", "Rome", "four stars", "gym", and "tennis" (or "hotel", "Roma", "quattro stelle", "palestra", and "tennis") over the Italian version of the Google on Jan. 6, 2008. Web pages returned by Google are reported in their original order in Table 5. As remarked in Section 1, it can be observed how there exist possibly out-of-scope pages

which have been ranked as very relevant (a four-star hotel without tennis facilities located in Abano Terme, five hundred kilometers from Rome is on the top of the result set) while potentially interesting pages (like the Rome Hilton Cavalieri) are positioned at the end of the list.

In order to apply our ranking methodology and show how existing search engines could benefit from its application, we manually annotated each page using concepts *Accommodation*, *Destination*, *Accommodation rating* and *Activity* in the *travel.owl* ontology. Relations were specified by following a fair approach relying only on information contained in the Web page (fourth column, keywords/concepts numbered progressively from one to five). Constrained page relevance score, is reported in column five. From column two, it can be easily observed that the ranking is significantly improved. For example, the first four entries now refer to Web pages which completely satisfy user query; entries five and six refer to hotels with all the requested characteristics, located in the vicinity of Rome and providing transfers to it. Another interesting example is provided by the fifth entry (now fifteenth), whose ranking was boosted through hidden text Web spam [12]. Even if Semantic Web will require the development of ad-hoc techniques for "semantic" anti-spam [31], our solution proved to be able to cope with the presence of malicious information in today's Web pages.

## 8 CONCLUSION

Next generation Web architecture represented by Semantic Web will provide adequate instruments for improving search strategies and enhance the probability of seeing user query satisfied without requiring tiresome manual refinement. However, actual methods for ranking returned result set will have to be adjusted to fully exploit additional contents characterized by semantic annotations

TABLE 3  
STATISTICS RELATED TO THE EXPERIMENTAL DATABASE USED FOR EVALUATING PERFORMANCE (TIME COMPLEXITY)

	$c_i (j=1)$	$c_i (j=2)$	$c_i (j=3)$	$c_i (j=4)$	$c_i (j=5)$	
$c_i (i=1)$	-	25131* 16811± 8305/8373/8482 <sup>□</sup>	25231* 20201± 10132/9988/9918/10013/10112 <sup>□</sup>	25097* 16793± 8416/8456/8346 <sup>□</sup>	25153* 0± 0 <sup>□</sup>	50186 <sup>#</sup>
$c_i (i=2)$	-	-	25118* 18855± 9251/9436/9412/9324 <sup>□</sup>	25119* 18940± 9519/9464/9475/9566 <sup>□</sup>	25062* 16750± 8351/8320/8256 <sup>□</sup>	50146 <sup>#</sup>
$c_i (i=3)$	-	-	-	25143* 16764± 8383/8297/8353 <sup>□</sup>	25081* 0± 0 <sup>□</sup>	49978 <sup>#</sup>
$c_i (i=4)$	-	-	-	-	25041* 0± 0 <sup>□</sup>	50181 <sup>#</sup>
$c_i (i=5)$	-	-	-	-	-	49990 <sup>#</sup>

TABLE 4  
EXPERIMENTAL RESULTS OF OUR APPROACH COMPARED WITH THOSE OF [16] (TIME COMPLEXITY)

Query concepts/ keywords	Relations among concepts	Sub-graphs processed [16]	CRKs processed [16]	Gen. CRKs/ Subm. to DB and inters. [16] (ms)	Total delay of [16] (ms)	Average num. of PFs ( $l=1/2/3/4$ )	Query the DB and gen. sub-graphs/ Proc. PFs (ms)	Total delay [our] (ms)
1	0	0	0	-/-	-	-/-/-/-	-/-	-
2	1	2	3	0.87/0.81	<b>1.68</b>	1/-/-/-	7.10/2.45	<b>9.55</b>
3	2	4	15	7.45/175.67	<b>183.12</b>	2/1/-/-	31.32/47.63	<b>78.95</b>
4	4	16	383	41.27/1251.37	<b>1292.64</b>	4/5/2/-	62.44/156.74	<b>219.18</b>
5	7	128	38399	632.41/2815.71	<b>3448.12</b>	7/12/11/9	84.69/688.16	<b>772.85</b>

Pages: # with at least one pair  $(c_i, k_j)$ , \* with both pairs  $(c_i, k_j)$  and  $(c_j, k_i)$ , ± including also at least one of the relations  $\bar{r}_{ij}^d$ , □ with specific relation  $\bar{r}_{ij}^d$ .

TABLE 5  
ACCURACY OF OUR RANKING ALGORITHM OVER THE FIRST TWENTY ENTRIES OF A RESULT SET GENERATED BY GOOGLE

Google ranking	Our ranking	Web page highlights, URL and content summary (some of the reasons for inclusion in the Google result set and guidelines used for identifying relations).	Rel. in page sub-graph	Constrained relev. score
1	13	<a href="#">Hotel 4 stelle Abano Terme - Hotel quattro stelle Abano Terme</a> <a href="http://www.abanoprenotazioni.it/hotel-4-stelle/Abano-terme.asp">www.abanoprenotazioni.it/hotel-4-stelle/Abano-terme.asp</a> Web page of a four-star hotel named "Hotel Terme Roma", with gym, in Abano Terme, five hundred kilometers from Rome. Another hotel in the Web page has gym and tennis courts.	$\delta_{13} = 1$ $\delta_{14} = 1$	$P(Q.p.4) = 0.000000$ $P(Q.p.3) = 0.000000$ $P(Q.p.2) = 0.333333$ $P(Q.p.1) = 0.428571$
2	14	<a href="#">Hotel 4 stelle Abano Terme - Hotel quattro stelle Abano Terme</a> <a href="http://www.abanoprenotazioni.it/hotel-4-stelle/Abano-terme.asp">www.abanoprenotazioni.it/hotel-4-stelle/Abano-terme.asp</a> Web page of a four-star hotel named "Hotel Terme Roma", with gym, in Abano Terme, five hundred kilometers from Rome. Another hotel in the Web page has gym and tennis courts.	$\delta_{13} = 1$ $\delta_{14} = 1$	$P(Q.p.4) = 0.000000$ $P(Q.p.3) = 0.000000$ $P(Q.p.2) = 0.333333$ $P(Q.p.1) = 0.428571$
3	11	<a href="#">Hotel 4 stelle Ischia - Hotel quattro stelle Ischia</a> <a href="http://www.ischiaprenotazioni.it/hotel-4-stelle-ischia.asp">www.ischiaprenotazioni.it/hotel-4-stelle-ischia.asp</a> Four-star hotel located in Ischia, renowned isle in front of Naples. Hotel facilities include gym and tennis courts.	$\delta_{13} = 1$ $\delta_{14} = 1$ $\delta_{15} = 1$	$P(Q.p.4) = 0.000000$ $P(Q.p.3) = 0.166667$ $P(Q.p.2) = 0.333333$ $P(Q.p.1) = 0.547619$
4	4	<a href="#">Hotel Petra &amp; Residence - Hotel Roma</a> <a href="http://www.initalia.it/hotel/hotelpetresidence.htm">www.initalia.it/hotel/hotelpetresidence.htm</a> Four-star hotel located in a residential area south-west of Rome, with gym and tennis courts. Easy connections to the city centre through public transports.	$\delta_{12} = 1$ $\delta_{13} = 1$ $\delta_{14} = 1$ $\delta_{15} = 1$	$P(Q.p.4) = 0.041667$ $P(Q.p.3) = 0.108974$ $P(Q.p.2) = 0.255556$ $P(Q.p.1) = 0.511905$
5	15	<a href="#">Roma vacanze on line viaggi città arte alberghi hotels bed ...</a> <a href="http://www.eurovacanza.com/strutture_html/idDest/9_lin/ita">www.eurovacanza.com/strutture_html/idDest/9_lin/ita</a> Four-star hotel located in Rome. Web page describes other hotels with different accommodation ratings. Web page has hidden spam including generic "gym" and "tennis" keywords.	$\delta_{12} = 1$ $\delta_{13} = 1$	$P(Q.p.4) = 0.000000$ $P(Q.p.3) = 0.000000$ $P(Q.p.2) = 0.250000$ $P(Q.p.1) = 0.892857$
6	5	<a href="#">Grand Hotel Duca d'Este Tivoli Roma</a> <a href="http://annunci.repubblica.it/roma/turismo/hotel/grand-hotel-duca-d-este-tivoli-pj-798427.html">annunci.repubblica.it/roma/turismo/hotel/grand-hotel-duca-d-este-tivoli-pj-798427.html</a> Four-star hotel located in Tivoli, a small town approximately twenty kilometers from Rome. Hotel facilities include gym and tennis courts.	$\delta_{12} = 1$ $\delta_{13} = 1$ $\delta_{14} = 1$ $\delta_{15} = 1$	$P(Q.p.4) = 0.041667$ $P(Q.p.3) = 0.108974$ $P(Q.p.2) = 0.255556$ $P(Q.p.1) = 0.511905$
7	8	<a href="#">Alberghi per sportivi a SORU E MARI, hotel con palestra, piscina ...</a> <a href="http://www.superdossier.com/attrezzati_per_lo_sport/SARDEGNA/CAGLIARI/SORU_E_MARI/">www.superdossier.com/attrezzati_per_lo_sport/SARDEGNA/CAGLIARI/SORU_E_MARI/</a> Four-star hotel with gym and tennis courts located in Sardegna. In the Web page there are links allowing to search for hotels in all the main Italian cities (Rome, and others).	$\delta_{14} = 1$ $\delta_{15} = 1$	$P(Q.p.4) = 0.000000$ $P(Q.p.3) = 0.166667$ $P(Q.p.2) = 0.333333$ $P(Q.p.1) = 0.547619$
8	6	<a href="#">Piste sci trekking scuola hotel albergo quattro stelle Monte ...</a> <a href="http://www.hotelcristalloterminillo.it/sport.asp">www.hotelcristalloterminillo.it/sport.asp</a> Four-star hotel with gym and tennis courts on the Monte Terminillo, a renowned winter location a few kilometers from Rome.	$\delta_{12} = 1$ $\delta_{13} = 1$ $\delta_{14} = 1$ $\delta_{15} = 1$	$P(Q.p.4) = 0.041667$ $P(Q.p.3) = 0.108974$ $P(Q.p.2) = 0.255556$ $P(Q.p.1) = 0.511905$
9	1	<a href="#">CROWNE PLAZA ROMA - ST. PETER'S HOTEL SPA - Hotel invest Italiana</a> <a href="http://www.hotel-invest.com/index.asp?id=385">www.hotel-invest.com/index.asp?id=385</a> Four-star hotel located in the city centre of Rome, with gym and tennis courts. The hotel is close to major attractions and offers a shuttle bus to main locations in the city.	$\delta_{12} = 3$ $\delta_{13} = 1$ $\delta_{14} = 1$ $\delta_{15} = 1$	$P(Q.p.4) = 0.125000$ $P(Q.p.3) = 0.224359$ $P(Q.p.2) = 0.366667$ $P(Q.p.1) = 0.583333$
10	12	<a href="#">Pagine gialle aziendali Viaggi turismo e tempo libero per l'import ...</a> <a href="http://viaggi-affari.europages.it/epq/dmc/1-it/did-21/hc-21510/Hotels_Italia.html">viaggi-affari.europages.it/epq/dmc/1-it/did-21/hc-21510/Hotels_Italia.html</a> Web page describing a four-star hotel in Rome, as well as other hotels located in different Italian cities with gym or tennis courts characterized by various accommodation ratings.	$\delta_{12} = 2$ $\delta_{13} = 1$	$P(Q.p.4) = 0.000000$ $P(Q.p.3) = 0.000000$ $P(Q.p.2) = 0.500000$ $P(Q.p.1) = 0.928571$
11	9	<a href="#">Hotel St. Gregory Park - Albergo quattro stelle San Giuliano Mare ...</a> <a href="http://www.abcfiere.com/dettaglio_hotel_st_gregory_park_rimini_219-1.php">www.abcfiere.com/dettaglio_hotel_st_gregory_park_rimini_219-1.php</a> Four-star hotel with gym and tennis courts located in Rimini (four hundred kilometers from Rome). Web page advertises several fashion events in Rome.	$\delta_{13} = 1$ $\delta_{14} = 1$ $\delta_{15} = 1$	$P(Q.p.4) = 0.000000$ $P(Q.p.3) = 0.166667$ $P(Q.p.2) = 0.333333$ $P(Q.p.1) = 0.547619$
12	3	<a href="#">Alberghi 4 stelle Roma - Hotel Aldobrandeschi - hotel Roma quattro ...</a> <a href="http://www.hotelaldobrandeschi.it/travel/it/alberghi_4_stelle_roma.htm">www.hotelaldobrandeschi.it/travel/it/alberghi_4_stelle_roma.htm</a> Four-star hotel in Rome with easy connections to the city and arranged rates for a gym in the vicinity. Hotel is close to a sporting complex with tennis courts.	$\delta_{12} = 2$ $\delta_{13} = 1$ $\delta_{14} = 1$ $\delta_{15} = 1$	$P(Q.p.4) = 0.083333$ $P(Q.p.3) = 0.166667$ $P(Q.p.2) = 0.311111$ $P(Q.p.1) = 0.547619$
13	2	<a href="#">hotel lusso roma</a> <a href="http://www.rome-luxury-hotel.com/hotel-lusso-roma.htm">www.rome-luxury-hotel.com/hotel-lusso-roma.htm</a> Four-star hotel located few minutes from the historical centre of Rome, with free shuttle bus to downtown. Hotel facilities include gym and tennis court.	$\delta_{12} = 2$ $\delta_{13} = 1$ $\delta_{14} = 1$ $\delta_{15} = 1$	$P(Q.p.4) = 0.083333$ $P(Q.p.3) = 0.166667$ $P(Q.p.2) = 0.311111$ $P(Q.p.1) = 0.547619$
14	10	<a href="#">Grand Hotel Palazzo della Fonte Fiuggi</a> <a href="http://www.iperhotel.com/index.cfm?Fuseaction=Hotel.InfoHotel&amp;IdHotel=325">www.iperhotel.com/index.cfm?Fuseaction=Hotel.InfoHotel&amp;IdHotel=325</a> Five-star hotel located in Fiuggi, a hundred kilometers from Rome, with gym and tennis courts (with instructors). Free shuttle bus to Rome city centre.	$\delta_{12} = 2$ $\delta_{14} = 1$ $\delta_{15} = 2$	$P(Q.p.4) = 0.000000$ $P(Q.p.3) = 0.166667$ $P(Q.p.2) = 0.393939$ $P(Q.p.1) = 0.833333$
15	18	<a href="#">GOLF TOSCANA - HOTEL - GOLF TOSCANA - CAMPI DA GOLF ...</a> <a href="http://www.hotelbenessere.it/golf-toscana.htm">www.hotelbenessere.it/golf-toscana.htm</a> Golf clubs in Toscana with gym and tennis courts. Web page provide links to hotels in the vicinity and driving directions from the main Italian cities (Rome, and others).	$\delta_{14} = 1$ $\delta_{15} = 1$	$P(Q.p.4) = 0.000000$ $P(Q.p.3) = 0.000000$ $P(Q.p.2) = 0.166667$ $P(Q.p.1) = 0.476191$
16	17	<a href="#">Hotel Cervia, Albergo Cervia, Hotels Cervia, Alberghi Cervia ...</a> <a href="http://it-hotel.7mates.com/italia/emilia-romagna/cervia.htm">it-hotel.7mates.com/italia/emilia-romagna/cervia.htm</a> Three-star hotel with gym named "Hotel Roma" located in Cervia, on the Riviera Romagnola. Web page describes other hotels (three and four stars) in Cervia with tennis courts.	$\delta_{14} = 1$ $\delta_{15} = 1$	$P(Q.p.4) = 0.000000$ $P(Q.p.3) = 0.000000$ $P(Q.p.2) = 0.166667$ $P(Q.p.1) = 0.476191$
17	19	<a href="#">Siria - Viaggi, valigie e salute</a> <a href="http://viaggiare.mevisi.it/articoli_viaggiare/2007011817361517436258.lasso">viaggiare.mevisi.it/articoli_viaggiare/2007011817361517436258.lasso</a> Travel to Siria departing from Rome and organized by a travel agency in Rome. Accommodation in a five-star hotel with gym and tennis courts (or in a four star hotel).	$\delta_{14} = 1$ $\delta_{15} = 1$	$P(Q.p.4) = 0.000000$ $P(Q.p.3) = 0.000000$ $P(Q.p.2) = 0.166667$ $P(Q.p.1) = 0.476191$
18	20	<a href="#">Mauritius, MAURITIUS - BLUE LAGOON BEACH 3 Stelle da 1viaggi ...</a> <a href="http://www.1viaggi.com/pacchetto/tour_fuga_nel-60.asp">www.1viaggi.com/pacchetto/tour_fuga_nel-60.asp</a> Travel to the Mauritius with accommodation in a three-star hotel with gym and tennis courts. Flight departing from Rome FCO airport.	$\delta_{14} = 1$ $\delta_{15} = 1$	$P(Q.p.4) = 0.000000$ $P(Q.p.3) = 0.000000$ $P(Q.p.2) = 0.166667$ $P(Q.p.1) = 0.476191$
19	7	<a href="#">Offerte - Hilton Cavalieri - Rome - Roma hotel - Italia - hotels.com</a> <a href="http://www.hotel.it/albergo-italia/albergo-roma/hilton-cavalieri-rome/">www.hotel.it/albergo-italia/albergo-roma/hilton-cavalieri-rome/</a> Five-star hotel located in the city centre of Rome, with free shuttle service to main city locations. Hotel facilities include gym and tennis courts.	$\delta_{12} = 2$ $\delta_{14} = 2$ $\delta_{15} = 2$	$P(Q.p.4) = 0.000000$ $P(Q.p.3) = 0.333333$ $P(Q.p.2) = 0.560606$ $P(Q.p.1) = 0.880952$
20	16	<a href="#">Roma vacanze on line viaggi città arte alberghi hotels bed ...</a> <a href="http://www.eurovacanza.com/strutture_html/idProvincia/94_lin/ita">www.eurovacanza.com/strutture_html/idProvincia/94_lin/ita</a> Four-star hotel located in Rome. Web page describes other hotels with different accommodation ratings. Web page has hidden spam including generic "gym" and "tennis" keywords.	$\delta_{12} = 1$ $\delta_{13} = 1$	$P(Q.p.4) = 0.000000$ $P(Q.p.3) = 0.000000$ $P(Q.p.2) = 0.250000$ $P(Q.p.1) = 0.892857$

including ontology-based concepts and relations. Several ranking algorithms for Semantic Web exploiting relation-based metadata have been proposed. Nevertheless, they mainly use page relevance criterions based on information that has to be derived from the whole knowledge base making their application often unfeasible in huge

semantic environments. In this work we propose a novel ranking strategy which is capable of providing a relevance score for a Web page into an annotated result set by simply considering user query, page annotation and underlying ontology. Page relevance is measured through a probability aware approach that relies on several graph-



based representations of the involved entities. By neglecting the contribution of the remaining annotated resources, a reduction in the cost of the query answering phase could be expected. Despite the promising results in terms of both time complexity and accuracy, further efforts will be requested to foster scalability into future Semantic Web repositories based on multiple ontologies, characterized by billions of pages, and possibly altered through next generation "semantic" spam techniques.

## REFERENCES

- [1] B. Aleman-Meza, C. Halaschek, I. Arpinar and A. Sheth, "A. Context-Aware Semantic Association Ranking," *Proc. Int'l Workshop on Semantic Web and Databases*, pp. 33--50, 2003.
- [2] K. Anyanwu, A. Maduko and A. Sheth, "SemRank: Ranking Complex Relation Search Results on the Semantic Web," *Proc. 14th Int'l Conf. on World Wide Web*, pp. 117-127, 2005.
- [3] R. Baeza-Yates, L. Calderón-Benavides and C. González-Caro, "The Intention Behind Web Queries," *Proc. 13th Int'l Conf. on String Proc. and Inf. Retr.*, pp. 98-109, 2006.
- [4] T. Berners-Lee and M. Fischetti, *Weaving the Web*. Harper Audio, New York, USA, 1999.
- [5] T. Berners-Lee, J. Hendler and O. Lassila, "The Semantic Web," *Scientific American*, 2001.
- [6] S. Brin and L. Page, "The Anatomy of a Large-Scale Hypertextual Web Search Engine," *Proc. 7th Int'l Conf. on World Wide Web*, pp. 107-117, 1998.
- [7] S. Cohen, J. Mamou, Y. Kanza and Y. Sagiv, "XSearch: A Semantic Search Engine for XML," *Proc. Int'l Conf. on Very Large Data Bases*, pp. 45-56, 2003.
- [8] L. Ding, T. Finin, A. Joshi, R. Pan, R. Scott Cost, Y. Peng, P. Reddivari, V. Doshi and J. Sachs, "Swoogle: a Search and Metadata Engine for the Semantic Web," *Proc. 13th ACM Int'l Conf. on Inform. and Knowledge Management*, pp. 652-659, 2004.
- [9] L. Ding, T. Finin, A. Joshi, Y. Peng, R. Pan and P. Reddivari, "Search on the Semantic Web," *IEEE Computer*, vol. 38, no. 10, pp. 62-69, 2005.
- [10] L. Ding, P. Kolari, Z. Ding and S. Avancha, "Using Ontologies in the Semantic Web: A Survey," *Ontologies*, Springer, pp. 79-113, 2007.
- [11] R. Guha, R. McCool and E. Miller, "Semantic Search," *Proc. 12th Int'l Conf. on World Wide Web*, pp. 700-709, 2003.
- [12] Z. Gyongyi and H. Garcia-Molina, "Spam: It's not Just for Inboxes Anymore," *IEEE Computer*, vol. 38, no. 10, pp. 28-34, 2005.
- [13] C. Junghoo, H. Garcia-Molina and L. Page, "Efficient Crawling through URL Ordering," *Computer Networks and ISDN Systems*, vol. 30, no. 1, pp. 161-172, 1998.
- [14] S. Kapoor and H. Ramesh, "Algorithms for enumerating all spanning trees of undirected and weighted graphs," *SIAM Journal on Computing*, vol. 24, pp. 247-265, 1995.
- [15] Y. Lei, V. Uren and E. Motta, "SemSearch: A Search Engine for the Semantic Web," *Proc. Managing Knowledge in a World of Networks*, pp. 238-245, 2006.
- [16] Y. Li, Y. Wang and X. Huang, "A Relation-Based Search Engine in Semantic Web," *IEEE Transactions on Knowledge and Data Engineering*, vol. 19, no. 2, pp. 273-282, 2007.
- [17] L. Page, S. Brin, R. Motwani and T. Winograd, "The PageRank Citation Ranking: Bringing Order to the Web," *Stanford Digital Library Technologies Project*, 1998.
- [18] A. Pisharody and H.E. Michel, "Search Engine Technique Using Keyword Relations," *Proc. Conf. on Art. Int.*, pp. 300-306, 2005.
- [19] T. Priebe, C. Schlager and G. Pernul, "A Search Engine for RDF Metadata," *Proc. 15th Int'l Workshop on Database and Expert Systems Applications*, pp. 168-172, 2004.
- [20] Protégé, Stanford Medical Informatics, H. Knublauch, <http://protege.cim3.net/file/pub/ontologies/travel/>, 2002.
- [21] Resource Description Framework (RDF) Model and Syntax Specification, <http://www.w3.org/TR/rdf-primer>, 2004.
- [22] C. Rocha, D. Schwabe and M.P. Aragao, "A hybrid approach for searching in the semantic web," *Proc. 13th Int'l Conference on World Wide Web*, pp. 374-383, 2004.
- [23] A. Sheth, C. Bertram, D. Avant, B. Hammond, K. Kochut and Y. Warke, "Managing Semantic Content for the Web," *IEEE Internet Computing*, pp. 80-87, 2002.
- [24] A. Sheth, B. Aleman-Meza, I. B. Arpinar, C. Bertram, Y. Warke, C. Ramakrishanan, C. Halaschek, K. Anyanwu, D. Avant, F. S. Arpinar and K. Kochut, "Semantic Association Identification and Knowledge Discovery for National Security Applications," *Journal of Database Management*, vol. 16, no. 1, pp. 33-53, 2005.
- [25] N. Stojanovic, "An Explanation-based Ranking Approach for Ontology-based Querying," *Proc. 14th Int'l Workshop on Database and Expert Systems Applications*, pp. 167-175, 2003.
- [26] N. Stojanovic, R. Studer and L. Stojanovic, "An Approach for the Ranking of Query Results in the Semantic Web," *Proc. Int'l Semantic Web Conference*, pp. 500-516, 2003.
- [27] R. Sun, H. Cui, K. Li, M.Y. Kan and T.S. Chua, "Dependency Relation Matching for Answer Selection," *Proc. 28th Int. ACM Conf. on Res. and Dev. in Inf. Retr.*, pp. 651-652, 2005.
- [28] T. Tran, P. Cimiano, S. Rudolph and R. Studer, "Ontology-based Interpretation of Keywords for Semantic Search," *Proc. 6th Int'l Semantic Web Conf.*, pp. 523-536, 2007.
- [29] Web Ontology Language, [www.w3.org/2004/OWL/](http://www.w3.org/2004/OWL/), 2004.
- [30] B.Y. Wu and K.M. Chao, *Spanning Trees and Optimization Problems*. CRC Press, New York, USA, 2004.
- [31] H. Yang, I. King and M.R. Lyu, "DiffusionRank: a Possible Penicillin for Web Spamming," *Proc. 30th Int'l ACM Conf. on Res. and Dev. in Inf. Retr.*, pp. 431-438, 2007.
- [32] Y.J. Zhang and Z.Q. Liu, "Refining Web Search Engine Results Using Incremental Clustering," *International Journal of Intelligent Systems*, vol. 19, no. 1, pp. 191-199, 2004.



**Fabrizio Lamberti** received his degree in computer engineering and his PhD degree in software engineering from the Politecnico di Torino, Italy, in 2000 and 2005, respectively. He has published a number of technical papers in international journal and conferences in the areas of distributed computing and information retrieval. He has served as a reviewer and program committee member for several international conferences and journals.



**Andrea Sanna** graduated in electronic engineering in 1993, and received his PhD degree in computer engineering in 1997, both from Politecnico di Torino, Italy. He has authored and coauthored several papers in the areas of computer graphics, scientific visualization, and parallel and distributed computing. He serves as reviewer for a number of international conferences and journals.



**Claudio Demartini** received his degrees in computer engineering and PhD in software engineering from the Politecnico di Torino, Italy, in 1980 and 1987, respectively. He has published a number of technical papers in international journals and conferences in the areas of distributed computing. He is the Vice Dean of the Industrial Engineering and Management School at Politecnico di Torino.